

Open Research Online

The Open University's repository of research publications and other research outputs

From *Dendroeca blackburniæ* to *Dendroeca blackburniæ*: what's in a name?

Conference or Workshop Item

How to cite:

King, David (2013). From *Dendroeca blackburniæ* to *Dendroeca blackburniæ*: what's in a name? In: TDWG 2013, 27 Oct - 1 Nov 2013, Florence, Italy.

For guidance on citations see [FAQs](#).

© 2013 ViBRANT

Version: Version of Record

Link(s) to article on publisher's website:

<https://mbgserv18.mobot.org/ocs/index.php/tdwg/2013/paper/view/344>

Copyright and Moral Rights for the articles on this site are retained by the individual authors and/or other copyright owners. For more information on Open Research Online's data [policy](#) on reuse of materials please consult the policies page.

oro.open.ac.uk

From *Dendroeca blackburniæ* to *Dendroeca blackburniae*: what's in a name?

Developing a names-based architecture assumes you have good, clean names to work with. While this assumption generally holds true for modern born-digital literature, the process of digitising legacy literature can produce errors.

Therefore, when extending the names-based architecture back in time it is necessary to take into account these errors.



An example: *Dendroeca blackburniæ*.

Dendroeca blackburniæ.
Dendroeca blackburniæ,
D. blackburniæ

This text represents a common challenge for optical character recognition (OCR) because the words are not in a dictionary, so cannot be automatically verified as they are processed.

Three examples of the name from the same text with:

- a variety of formats,
- different sizes,
- stroke weights,
- directions to the face, and
- three different forms of the æ ligature.

The two ligatures, æ and œ, present an additional challenge. The OCR engine was set to expect modern English text, but the ligatures do not appear in modern English and so can **never** be recognised by the OCR process.

Fuzzy matching

Fuzzy matching can help correct the OCR rendering by finding similarities across the renderings and to the correct spelling.

Dendroeca occurs 59 times in the text, rendered as:

- Dendroeca 32 times,
- Dendreca 23 times, and
- once each for Bendrceca, Bendrwca, DendrcBca and Dendrosca.

The text also contains Dendroica, which occurs five times, correctly rendered by the OCR every time. Dendroeca and Dendroica will match fuzzily!

Therefore, we need collocation too.

Acknowledgements

This research uses the Biologia Centrali-Americana (BCA). PDFs and OCR renderings can be downloaded from the Biodiversity Heritage Library, www.biodiversitylibrary.org.

Thank you to Anna Weitzman and Chris Lyal of the INOTAXA project, www.inotaxa.org, for making their project's re-keyed texts of the BCA available for our research.

One final challenge: Try looking up *Dendroeca blackburniæ* in a modern taxonomic reference. But that's another project...

Collocation

This technique examines surrounding words to provide the context of use which helps disambiguate similar words.

Collocation can help with blackburniæ, which occurs six times in the text. The OCR recognises the word as:

- blackburnice four times,
- blackburniae once, and
- blackburnw once.

Collocation shows that blackburniæ – however it is spelt – follows what looks like a genus or genus abbreviation. This additional information allows us to target our name correction to plausible binomial combinations.

Read more about OCR post-processing

The ViBRANT project is building a corpus of marked up documents for research into OCR issues. This is freely available from git.scratchpads.eu/v/vibrantcorpus.git.

We are preparing papers from our work, covering the tools and workflows used in developing the corpus, and preliminary findings from analysis of the corpus.